

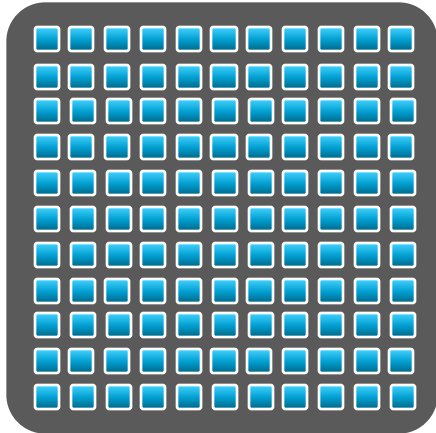
AI Networking - Requirements & Solutions

Kamran Naqvi

Chief Network Architect – EMEA, CSG

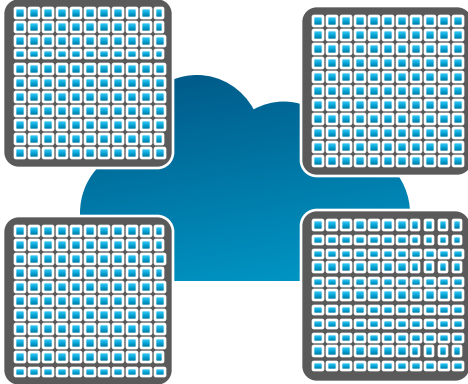
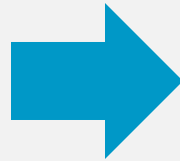
October 2025

AI ... A Very Very Large Distributed Computing System



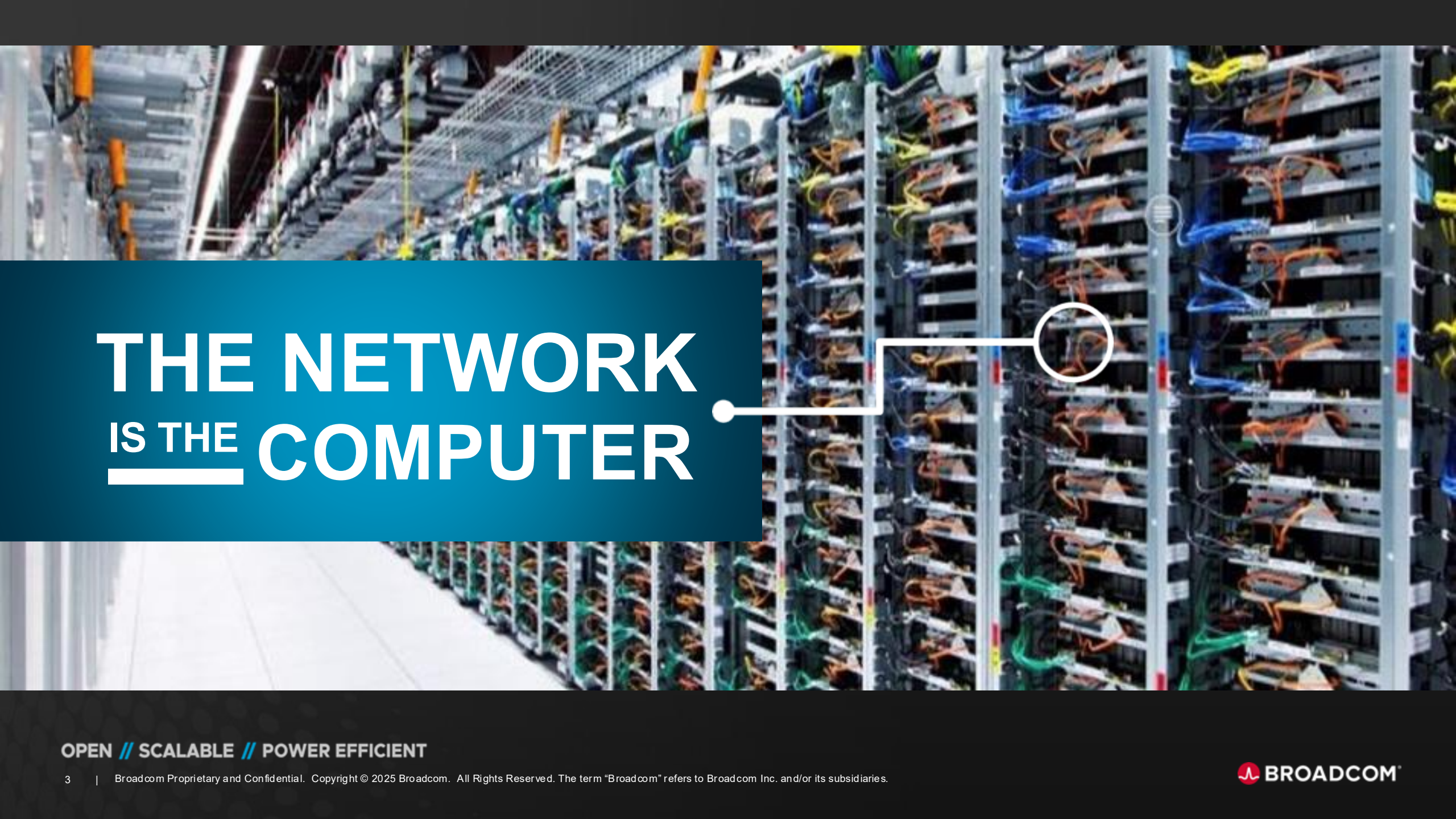
GPU
Optimized for
Parallel Tasks

Thousands of Cores



GPU Clusters
GPU Network for AI
Workloads

100K+ GPUs



THE NETWORK IS THE COMPUTER

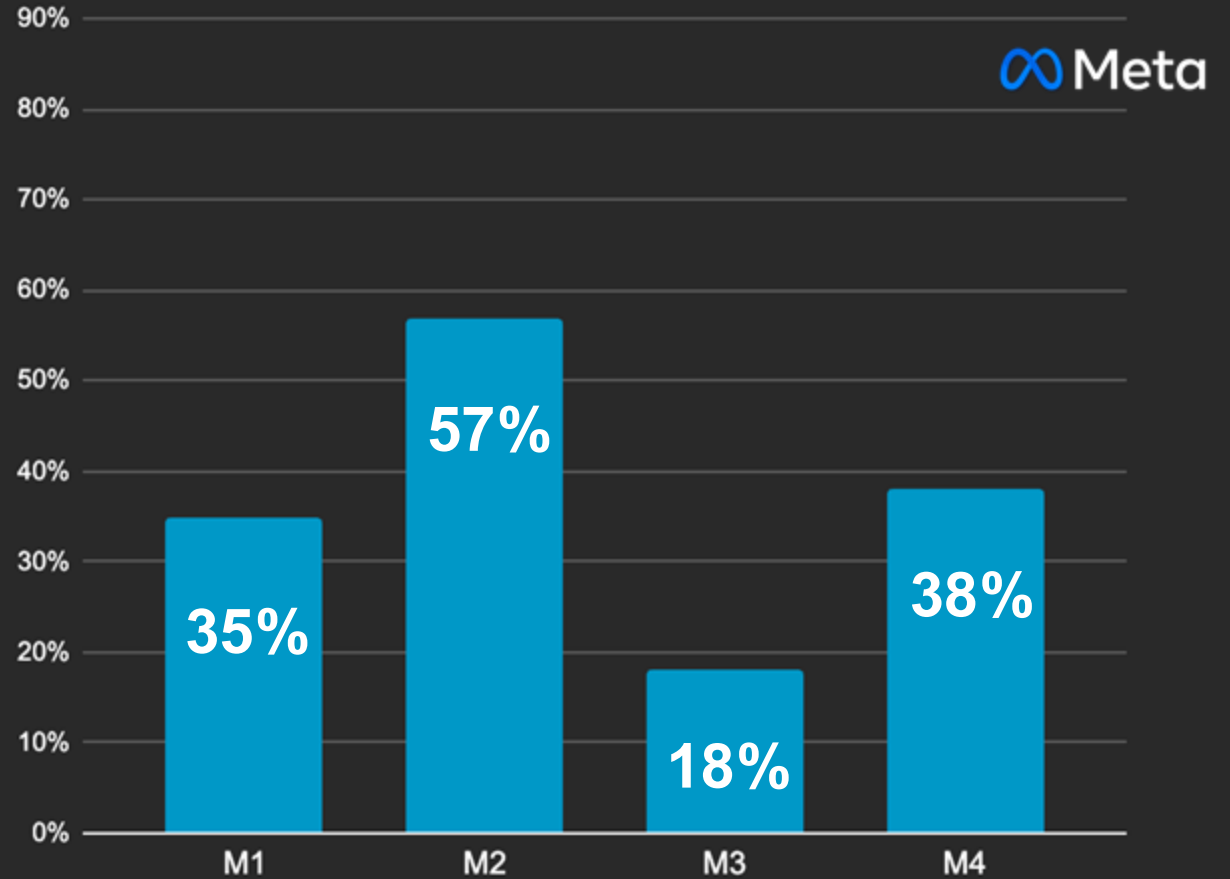
OPEN // SCALABLE // POWER EFFICIENT



Network I/O is Key for Recommendation Workloads. ”

OCP keynote by Alexis Bjorlin at 2022 OCP Global Summit

TIME SPENT IN NETWORKING

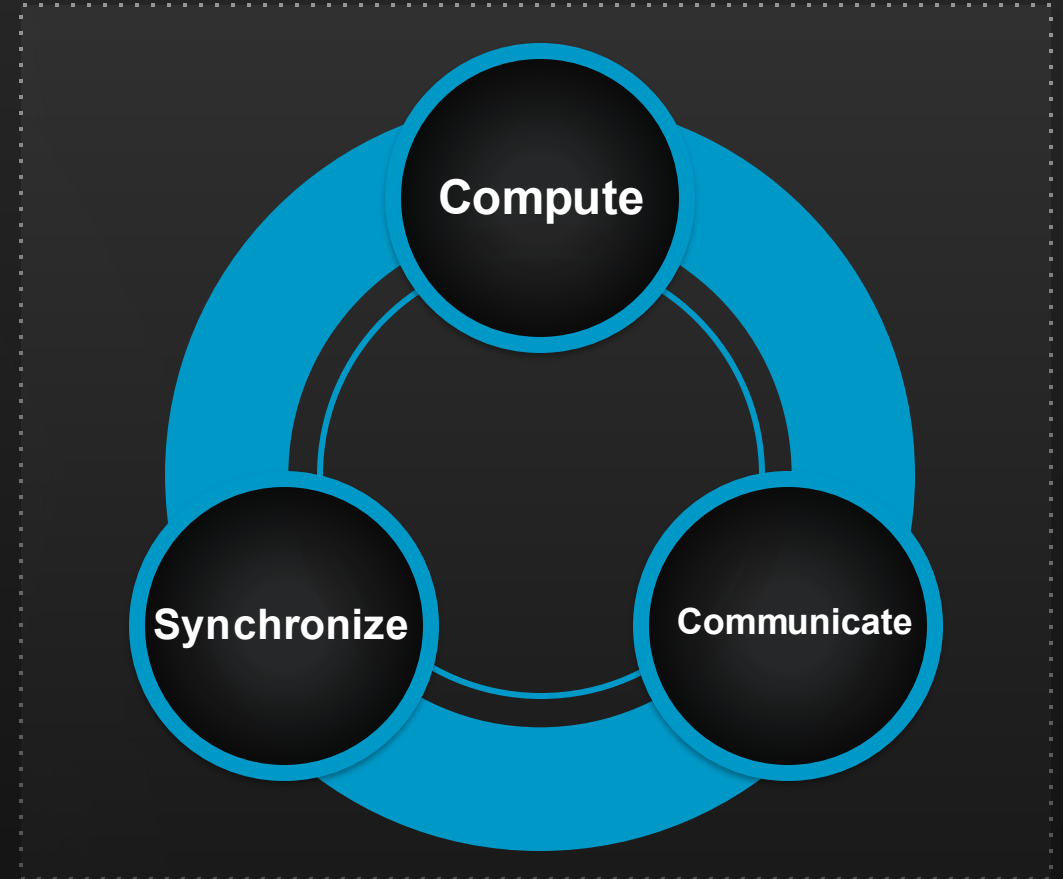


M# = ML model #

Ranking requires high injection & bisection bandwidth

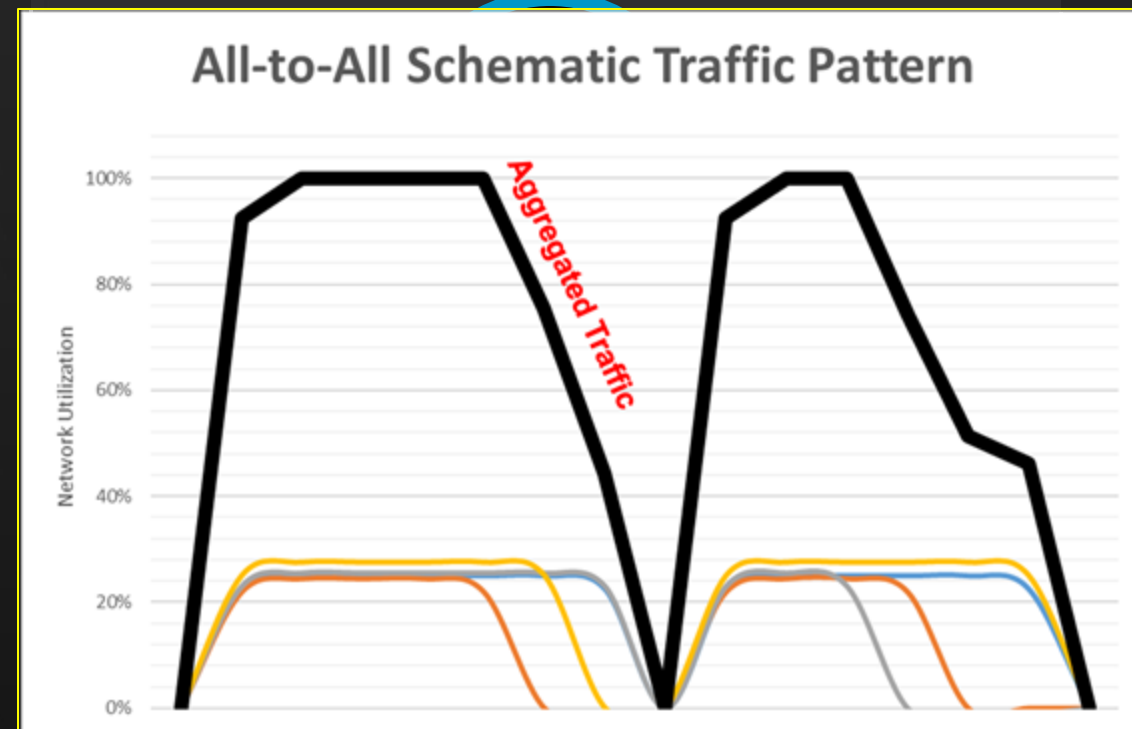
What Makes AI Networks Unique?

- High bandwidth
- Elephant flows
- Synchronized and bursty traffic
- RDMA dominant traffic
- Training jobs run for long periods of time (hours, days)
- Tail latency impacts job completion time significantly

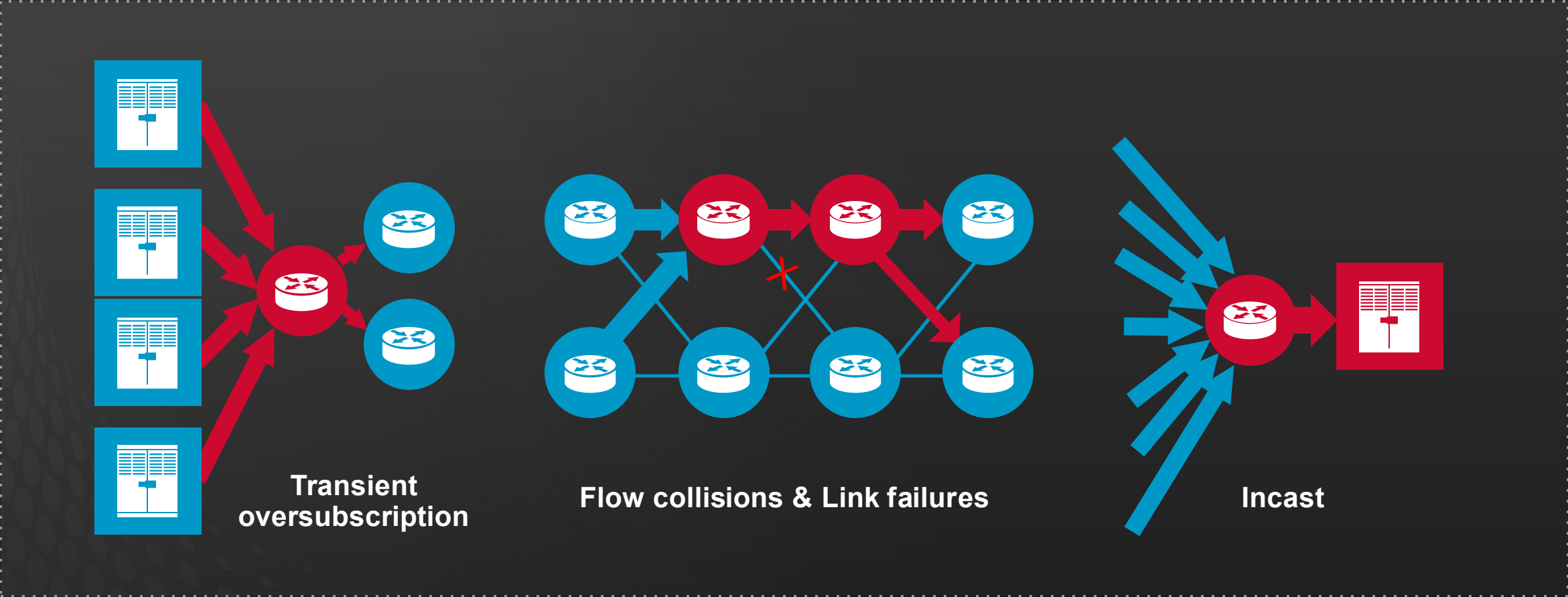


What Makes AI Networks Unique?

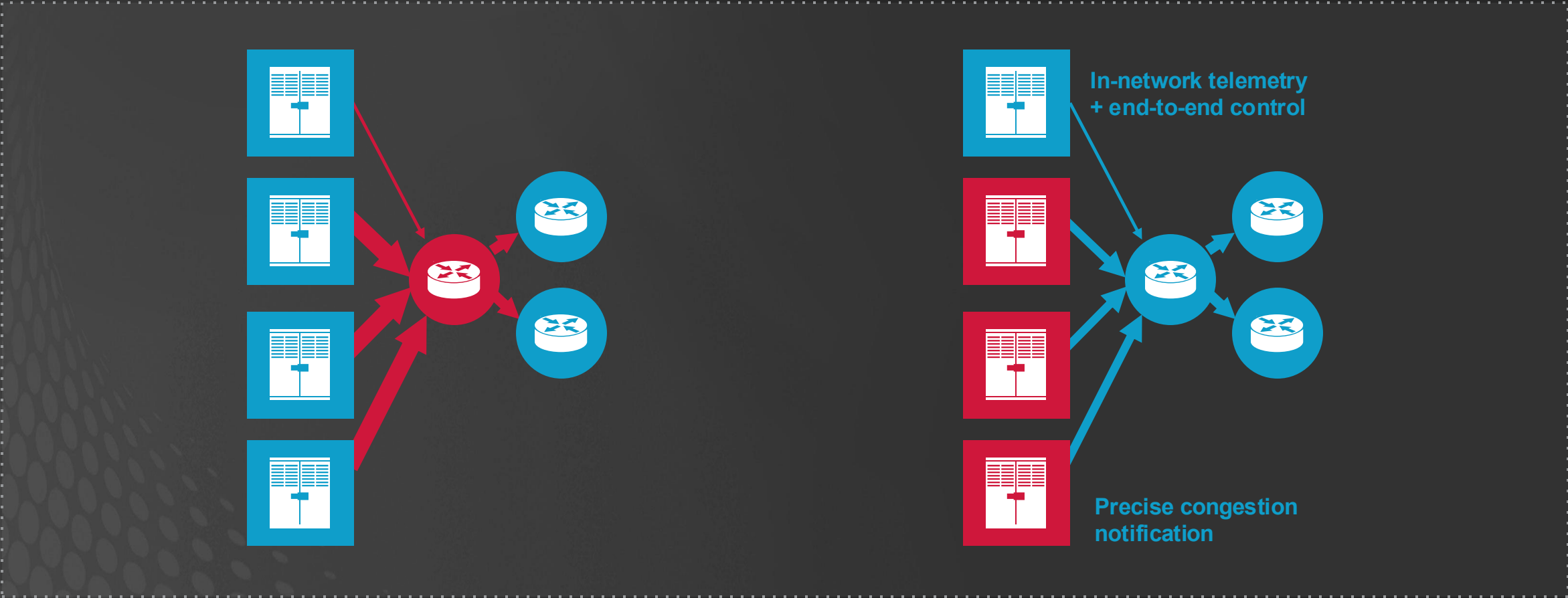
- High bandwidth
- Elephant flows
- Synchronized and bursty traffic
- RDMA dominant traffic
- Training jobs run for long periods of time (hours, days)
- Tail latency impacts job completion time significantly



“Time Spent in Networking” is Impacted By...



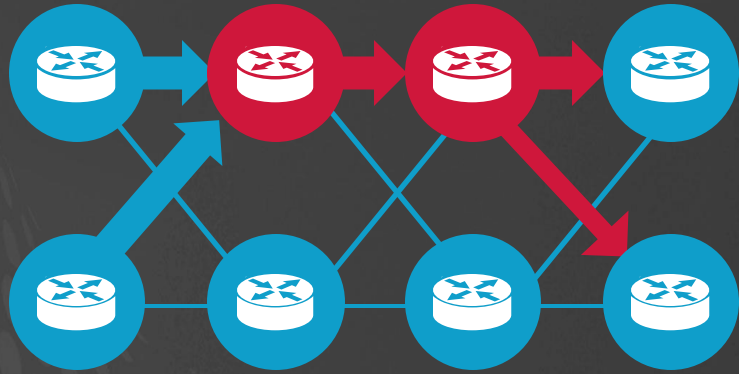
“Time Spent in Networking” is Improved By...



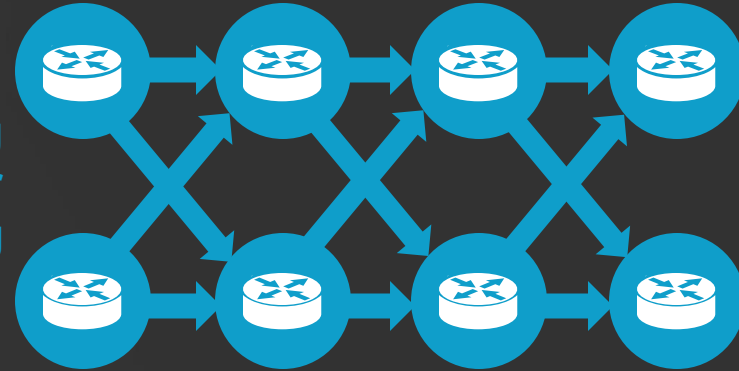
OPEN // SCALABLE // POWER EFFICIENT

| Copyright © 2025 Broadcom. All Rights Reserved. The term "Broadcom" refers to Broadcom Inc. and/or its subsidiaries.

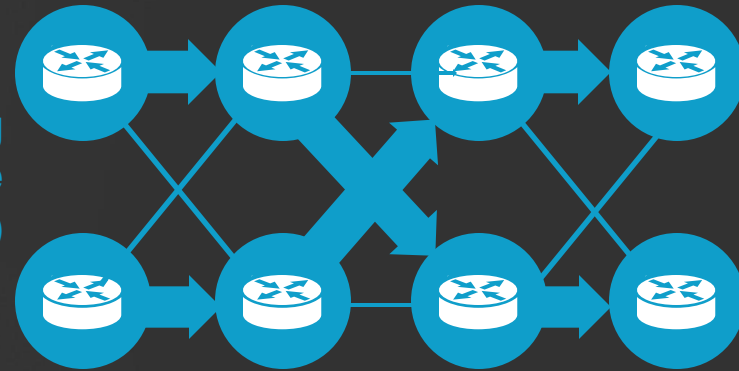
“Time Spent in Networking” is Improved By...



Packet spraying
with receiver
ordering



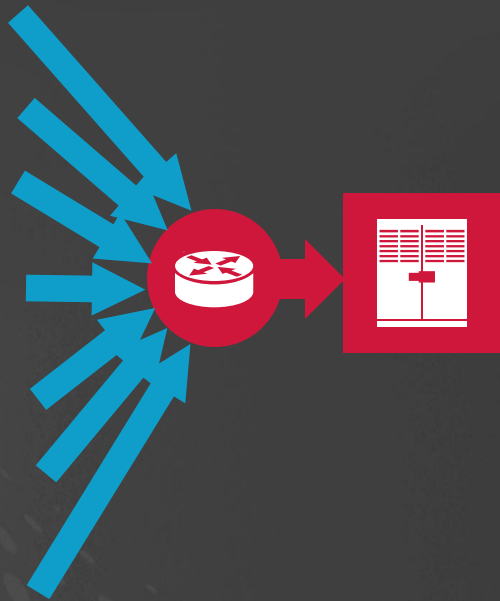
Cognitive routing
(load-aware
ECMP)



In case of link failure recovery should happen in hardware ... Zero Impact Failover (ZIF)

OPEN // SCALABLE // POWER EFFICIENT

“Time Spent in Networking” is Improved By...



Receiver-based credit control can pace senders accurately.

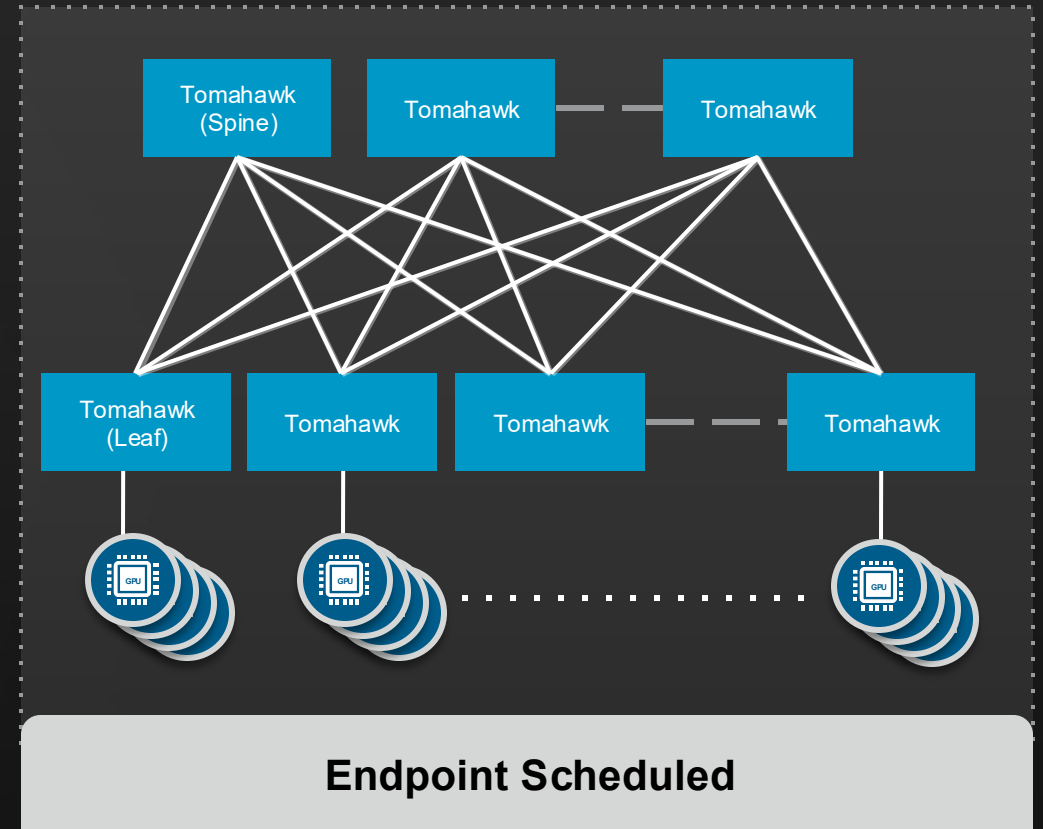
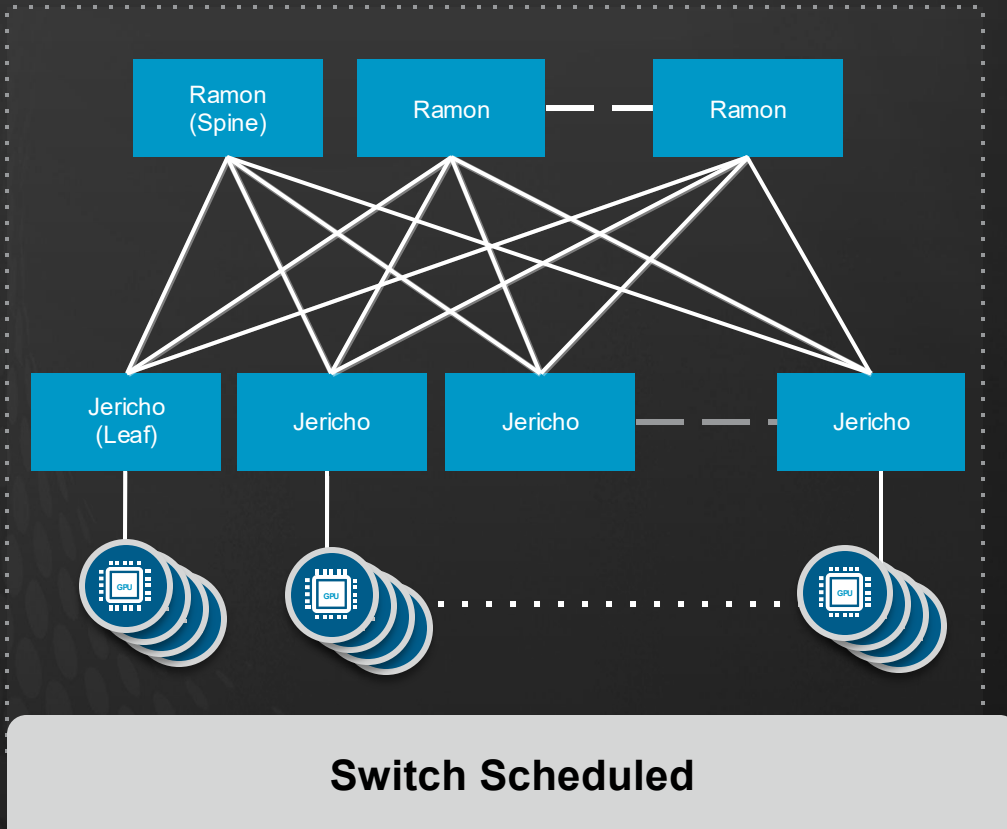
Credit control mechanism can exist on the switch or the endpoint

OPEN // SCALABLE // POWER EFFICIENT

| Copyright © 2025 Broadcom. All Rights Reserved. The term "Broadcom" refers to Broadcom Inc. and/or its subsidiaries.

 BROADCOM

Broadcom's AI Fabric Solutions



¹¹ OPEN // SCALABLE // POWER EFFICIENT

| Copyright © 2025 Broadcom. All Rights Reserved. The term "Broadcom" refers to Broadcom Inc. and/or its subsidiaries.

Tomahawk 6: Built for AI Scale



World's First 102.4 Tbps Switch Chip

- Double the bandwidth of any other Ethernet switch
- Built to power clusters with 1M+ XPU



Unmatched Performance & Power Efficiency

- Cognitive Routing 2.0, Deep Insight
- Advanced 3nm technology



Unrivaled Versatility

- Scale-up & scale-out, training & inference
- Works with any endpoint, including XPU scale-up interfaces



Industry-Leading SerDes and CPO

- Options for 512x200G PAM4, 1024x100G PAM4, CPO
- Ground-breaking SerDes and optics density

NOW SHIPPING



Tomahawk6-200G

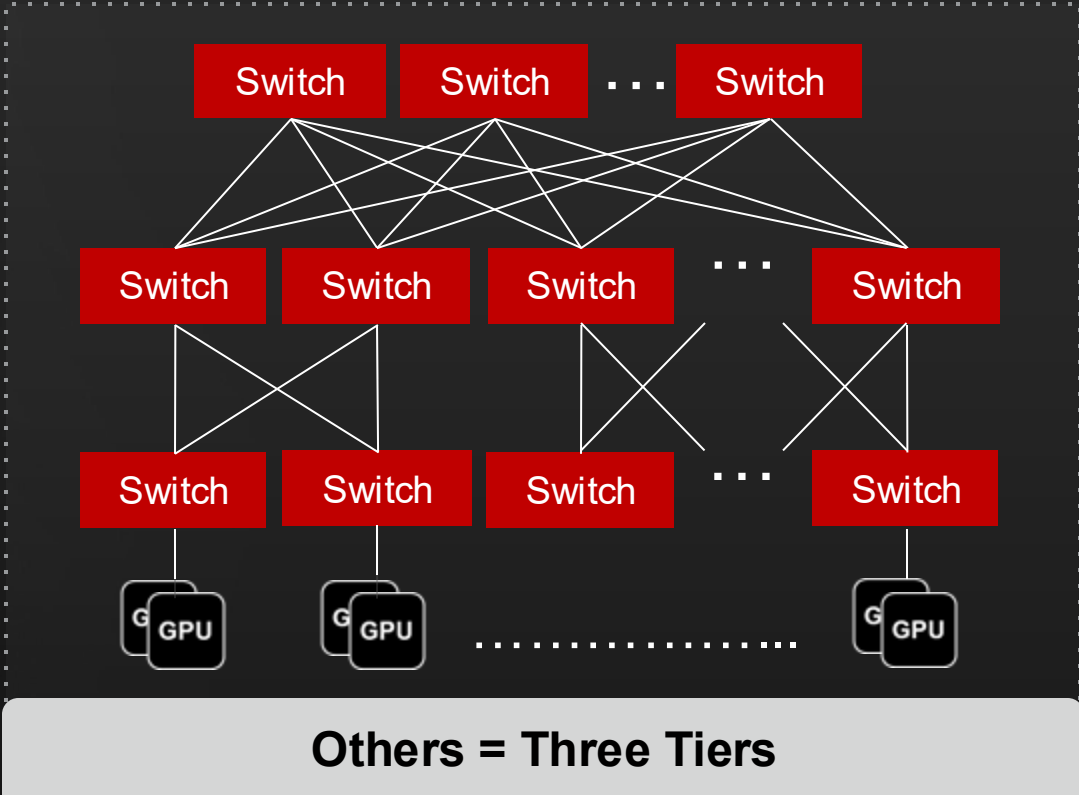
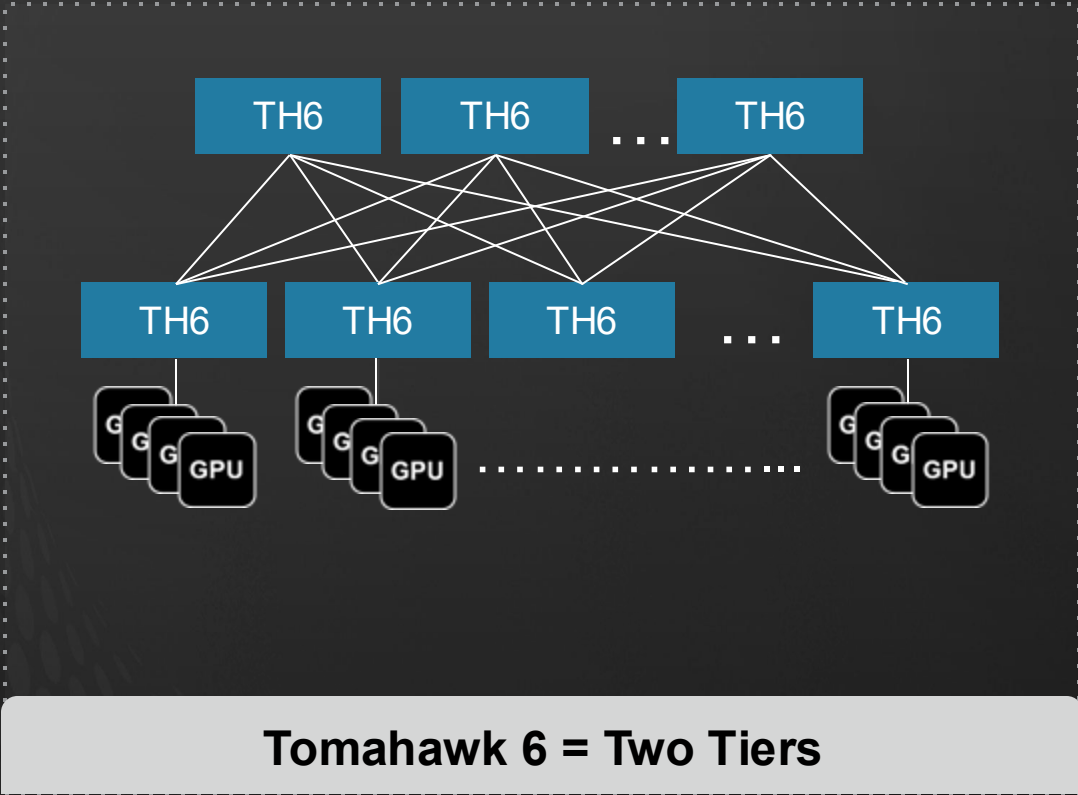


Tomahawk6-100G

Support for 1.6TbE Ports
Ultra Ethernet Compliant

OPEN // SCALABLE // POWER EFFICIENT

Tomahawk 6: Large Two-Tier Scale-Out



N = # rails

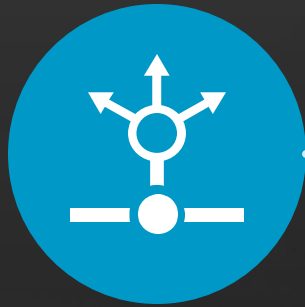
128K-GPU Cluster at 200GbE x N per Endpoint

OPEN // SCALABLE // POWER EFFICIENT

Benefits of Two-Tier vs. Three-Tier Networks

Two-Tier	Three-Tier
Fewer optics	67% more optics
Lower latency	67% higher latency (5 vs. 3 hops)
Higher reliability	>3x number of switches & 67% more optics
Higher performance	Difficult load balancing and congestion control, more link flaps
Lower power	2x power for the network

Broadcom Cognitive Routing 2.0



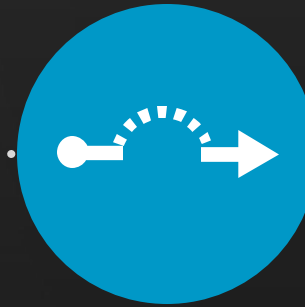
Global Load Balancing

Egress link selection based on global path congestion



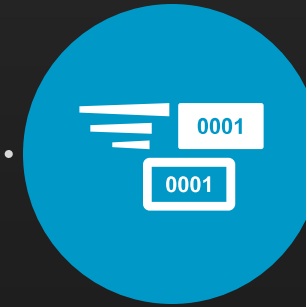
Reactive Path Rebalancing

Update egress links for active flows when congestion is detected



Fast Link Failover

Automatically steers traffic around failed links in <500ns



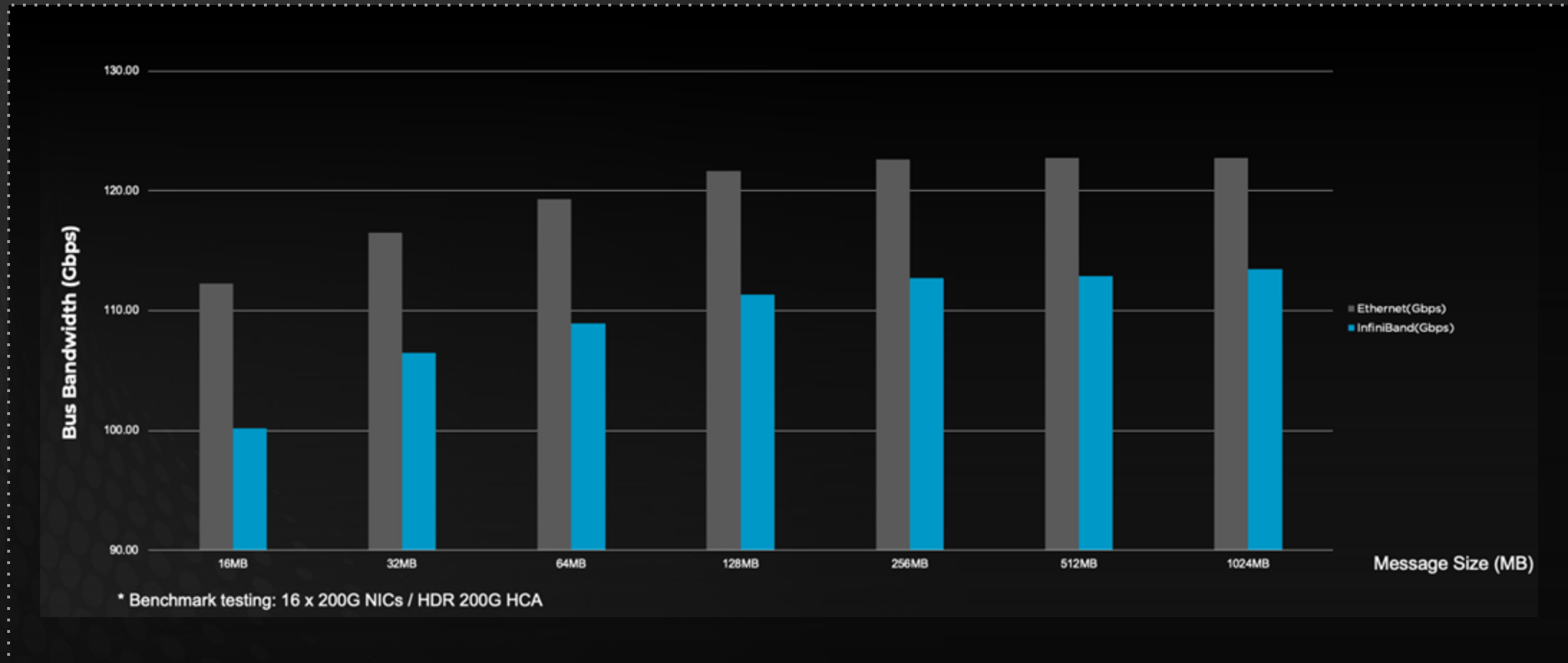
Drop Congestion Notification

For congested queues, send trimmed packet at high priority to destination

New for Cognitive Routing 2.0:

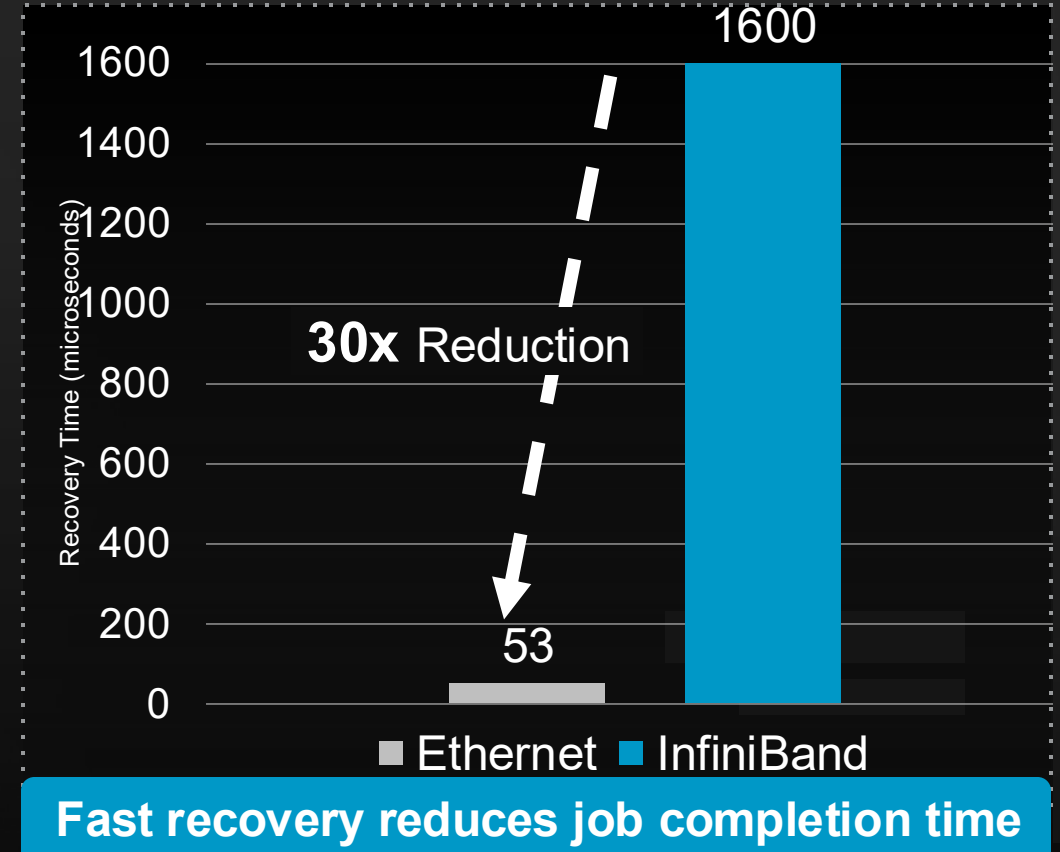
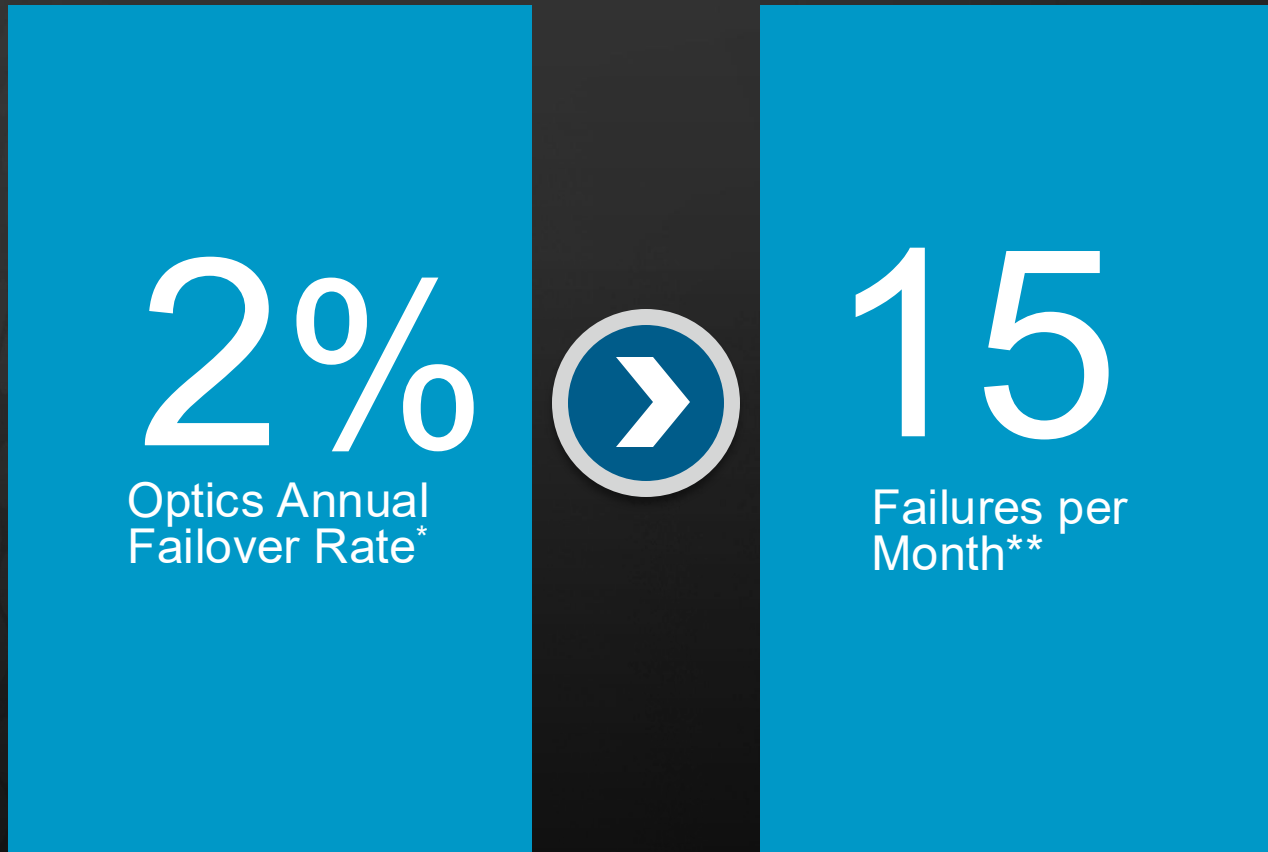
- Increased scale
- Faster and enhanced congestion signaling
- Enhanced path selection algorithms

Ethernet vs. InfiniBand: 10+% Improvement in Job Completion Time



OPEN // SCALABLE // POWER EFFICIENT

Ethernet Provides 30x Faster Failover than InfiniBand



* Typical industry failure rate. ** Assuming 4K node cluster using 9.2K optic modules

OPEN // SCALABLE // POWER EFFICIENT

All Hyperscalers: Ethernet AI fabric



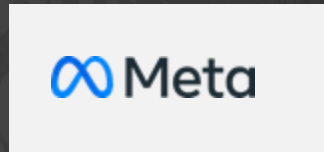
100,000+



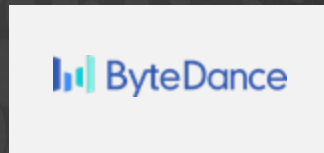
60,000+



30,000+



30,000+

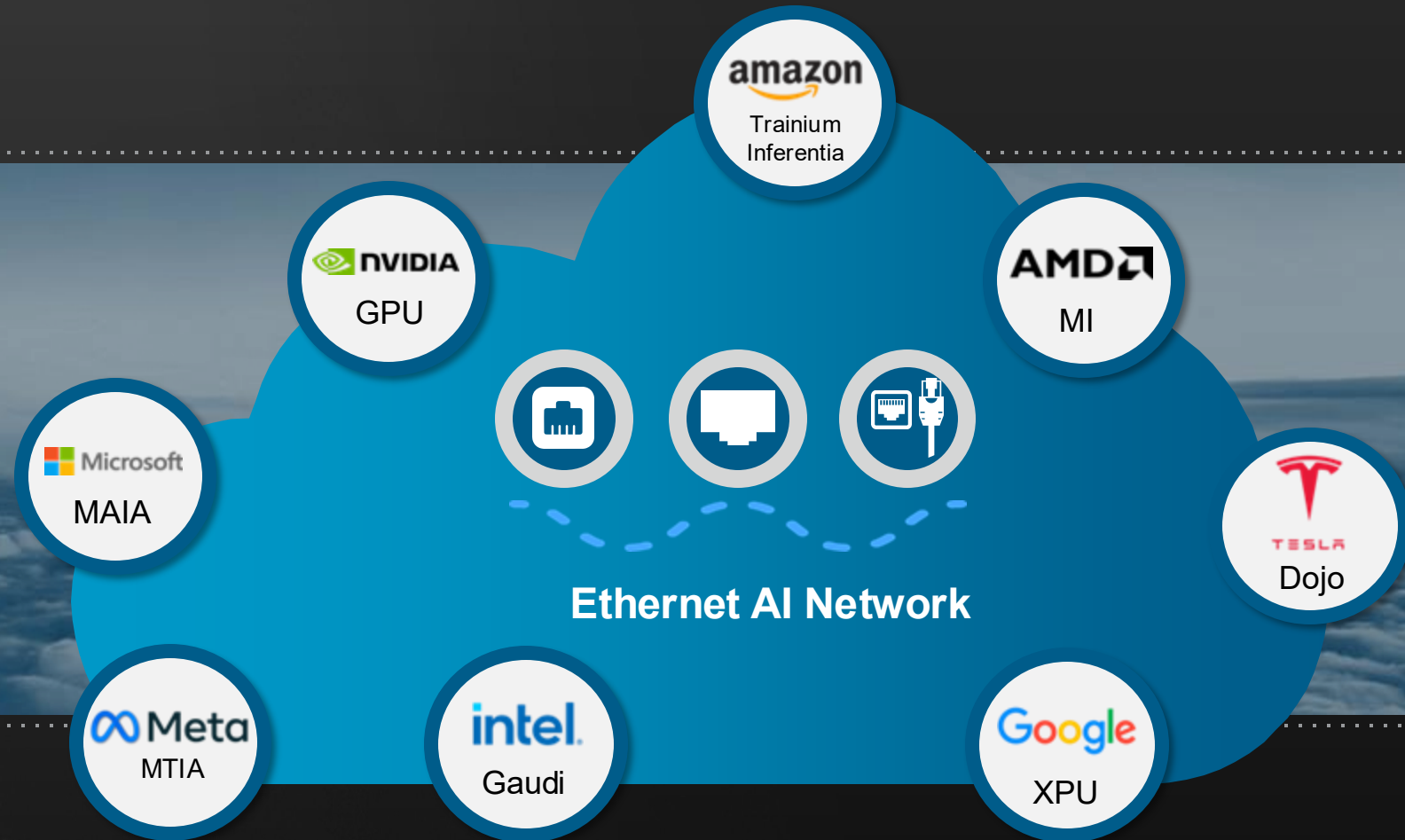


10,000+



OPEN // SCALABLE // POWER EFFICIENT

Ethernet Network for AI



¹⁹ OPEN // SCALABLE // POWER EFFICIENT

| Copyright © 2025 Broadcom. All Rights Reserved. The term "Broadcom" refers to Broadcom Inc. and/or its subsidiaries.



Thank You
